

# Year 4 — Numerical Linear Algebra

Based on lectures by Prof. Yuji Nakatsukasa

Notes taken by James Arthur

Michaelmas 2022

These notes are not endorsed by the lecturers, and I have modified them (often significantly) after lectures. They are nowhere near accurate representations of what was actually lectured, and in particular, all errors are almost surely mine (especially the typos!).

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Why do we care? . . . . .	2
1.2	Norms . . . . .	2

# 1 Introduction

The goal of this course is to solve and understand problems using (usually non-square) matrices. The course will cover,

1. SVD. Singular Value Decomposition, see Linear Algebra (Year 2).  $A = U\Sigma V^T$ .
2. Linear systems  $Ax = b$  and eigenvalue problems  $Ax = \lambda x$ . Least square problems,  $\min_x \|Ax - b\|_2$ .

There are going to be three classes of approaches,

- Direct Methods, classical approach. Not useful for very big matrices.
- Iterative solvers, work even if matrix sizes are larger.
- Randomised Algorithms, use some sense of randomisation in order to get an algorithm that works with high probability better than direct for very large matrices.

## 1.1 Why do we care?

When we want to solve a non-linear equation a linear system pops out. Motivation: minimising a function,  $\min_x f(x)$  where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , where we let  $n$  be large. One powerful way to find a minimum is to find a stationary point. So we find,

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} = 0$$

if we have a nice (convex) function, then if we have a stationary point, we have a minimiser. This boils down to, letting  $F = \nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ . We need to find zeros of this non-linear problem. Hence we now have a root-finding problem, and so we can use Newton's Method.

$$x_{\text{new}} = x_{\text{old}} - \mathcal{J}^{-1}F(x_{\text{old}})$$

and we see,

$$\mathcal{J}_{ij} = \frac{\partial F_i}{\partial x_j}$$

and then this is the hessian of  $f$ . This is then a linear system,

$$\mathcal{J}\Delta x = F(x_{\text{old}}).$$

**NB! Here  $A^* = \bar{A}$  (instead of  $A^* = \bar{A}^T$ )** We are going to stay under real matrices because there are only two cases where the difference matters. Further,  $m \geq n$  for most matrices in this course. For orthonormal matrices,

$$A^T A = I_n \quad AA^T \neq I_m$$

## 1.2 Norms

We need norms to quantify how big matrices are. These help us when approximating matrices. Given some  $\mathbf{x} \in \mathbb{R}^n$ , we have

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2}$$

and 1-norm,

$$\|x\|_1 = |x_1| + \cdots + |x_n|$$

and the p-norm,

$$\|x\|_p = (|x_1|^p + \cdots + |x_p|^p)^{\frac{1}{p}}$$

and finally the  $\infty$ -norm,

$$\|x\|_\infty = \max_i |x_i|$$

**Definition 1.1** (Norm). A norm satisfies three axioms,

- $\|\alpha x\| = |\alpha| \|x\|$
- $\|x\| \geq 0$  and  $\|x\| = 0 \iff x = 0$
- $\|x + y\| \leq \|x\| + \|y\|$

**Lemma 1.2** (Holder's Inequality). For  $p > q$ ,

$$\|x\|_p \leq \|x\|_q$$

**Definition 1.3** (Unitarily Invariant). If  $A$  is orthonormal, then  $\|Ax\|_2 = \|x\|_2$ .

**Lemma 1.4** (Cauchy-Schwartz). For any  $x, y \in \mathbb{R}^n$ ,

$$|x^T y| \leq \|x\|_2 \|y\|_2$$

*Proof.* For any scalar  $c$ ,  $\|x - cy\|_2^2 = \|x\|_2^2 - 2cc^T y + \|y\|_2^2$ . Now we complete the square and we get,

$$\begin{aligned} \|x - cy\|_2^2 &= \|x\|_2^2 - 2cc^T y + \|y\|_2^2 \\ &= \|y\|_2^2 \left( c - \frac{x^T y}{\|y\|_2} \right)^2 + \|x\|_2^2 - \frac{(x^T y)^2}{\|y\|_2^2} \end{aligned}$$

and so we now minimise by letting  $c = \frac{x^T y}{\|y\|_2^2}$  and so we get Cauchy Schwartz. □

Now for matrix norms. We have the  $p$ -norm,

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \frac{\|Ax\|_p}{\|x\|_p}.$$

The most important case is when  $p = 2$ , this is the spectrum norm. This is,

$$\|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$$

**Exercise.** Show,

$$\|A\|_1 = \text{maximum column sum}$$

$$\|A\|_\infty = \text{maximum row sum}$$

**Definition 1.5** (Frobenius Norm).

$$\|A\|_F = \sqrt{\sum_i \sum_j |A_{ij}|^2} = \sqrt{\mathbf{A}^T \mathbf{A}} = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})}$$

where

$$\mathbf{A} = \begin{pmatrix} A_{11} \\ \vdots \\ A_{1n} \\ A_{21} \dots \end{pmatrix}$$

**Definition 1.6** (Trace Norm).

$$\|A\|_* = \sum_{i=1}^{\min(m,n)} \sigma_i(A)$$

For most  $p$ -norms,

$$\|AB\|_p \leq \|A\|_p \|B\|_p$$

and for the Frobenius norm,

$$\begin{aligned} \|AB\|_F &\leq \|A\|_F \|B\|_F \\ &\leq \|A\|_2 \|B\|_F \end{aligned}$$

this is the subordinate property. Where this goes wrong is,

$$\|A\|_\infty = \max_{i,j} |A_{ij}|$$